# Research in
# Multi-Resource Aware Scheduling Algorithms*

## HPCMOD UGC, June 6th, 2000

William Leinberger, George Karypis, and Vipin Kumar
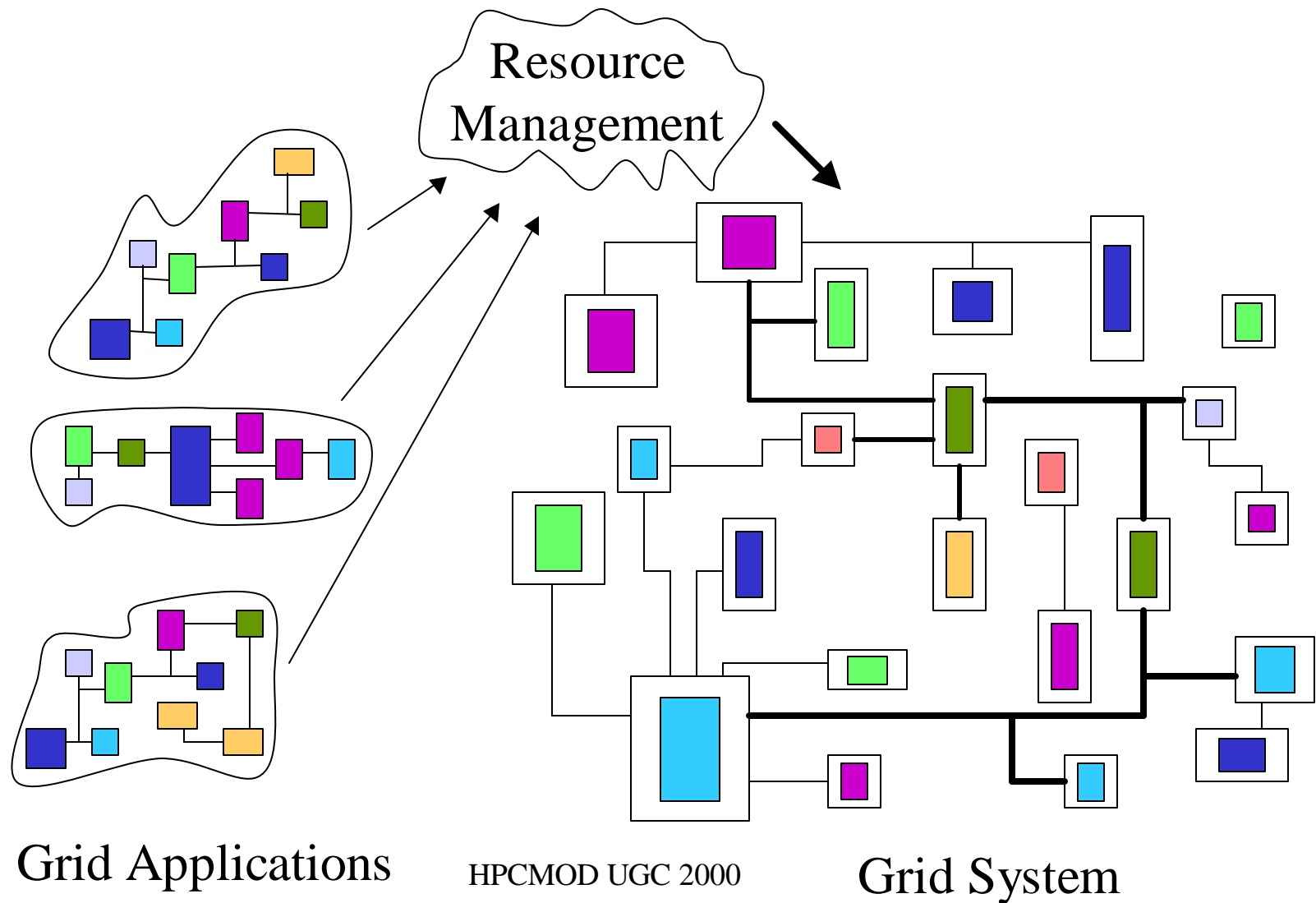
{leinberg, karypis, kumar} @ cs.umn.edu

Army High Performance Computing Research Center

University of Minnesota

# Scheduling Grid Applications: The Vision



Resource Management

Grid Applications

HPCMOD UGC 2000

Grid System

# Scheduling Grid Applications: Application-Centric Approach

- An application-specific scheduler is designed to:
    - Query the grid for resource candiates
    - Select the "best" resources, typically through advance reservation
    - Montior the resources during execution and adjust as necessary.

\+ Potentially better execution performance for the single application.

– Must be maintained along with the application

– Degrades to First-Come-First-Served resource allocation, resulting in low utilization of critical grid resources.

# Scheduling Grid Applications: System-Centric Approach

- System level scheduling:
  - Accepts resource requests from grid-applications
  - Allocates grid resources to all waiting jobs to make best use of all grid resources.

+ More efficient use of critical grid resources.

- Much harder to implement general system level scheduler due to scale of grid systems.

Our Research focuses on system level scheduling approaches!

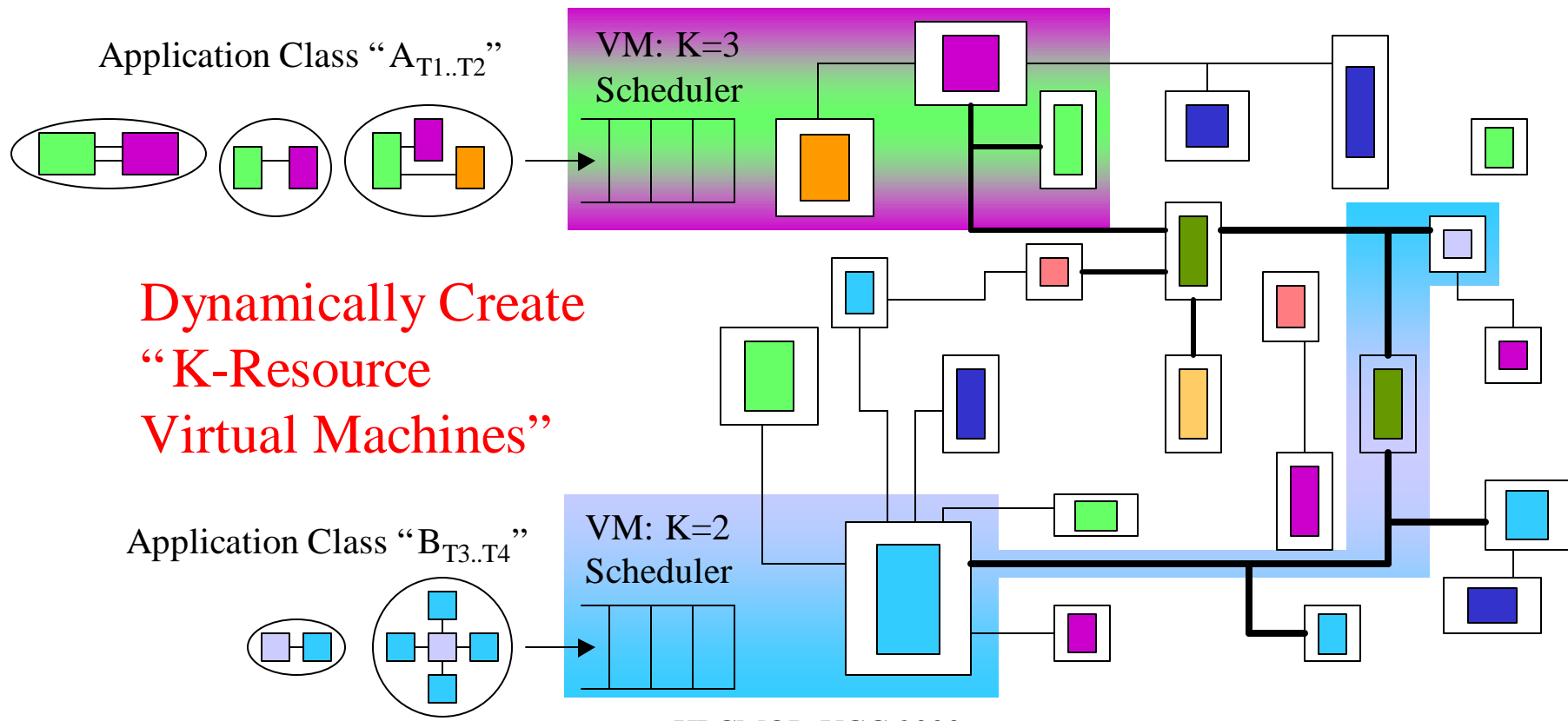# Scheduling Grid Applications: Assumptions

- Multiple instances of grid applications are submitted in common temporal space:
  - Parameter sweep studies
  - Real-time data collection and analysis

- Grid applications may be grouped into special resource configurations:
  - Functionally Related Resources: e.g. a data source, compute engine, and visualization interface.
  - Sparse Resources: e.g. peta-flop systems or data repositories, special purpose processor, high-speed data link
  - Atomic Resources: e.g. sensors, microscopes, etc.

We propose the use of dynamically created Virtual Machines!
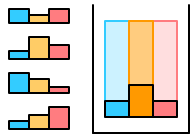
HPCMOD UGC 2000

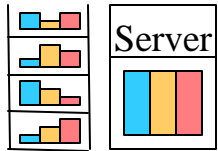# Scheduling Grid Applications: A Solution



Grid Applications

Grid System

Application Class "$A_{T1..T2}$"

VM: K=3
Scheduler

**Dynamically Create
"K-Resource
Virtual Machines"**

Application Class "$B_{T3..T4}$"

VM: K=2
Scheduler

HPCMOD UGC 2000

# Research Topics for Scheduling K-Resource Virtual Machines

Bin Packing
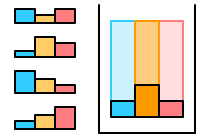
- The k-Capacity Bin Packing Problem
  - The Resource Balancing Heuristic

Scheduling

Server

- Scheduling on a Single K-Resource VM
  - Jobs are items, free resources are a bin

Multi-Server

- Load Balancing Across Multiple Near-Homogeneous K-Resource VMs
  - In case one VM isn't enough

Multi-Resource

- Dynamic Creation of K-Resource VMs in a Computational Grid
  - How and Who

HPCMOD UGC 2000
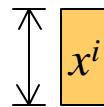
# K-Capacity Bin Packing: Notations

## Single Capacity          *k*-Capacity

**Item List L:**     $\{x^1, x^2, ....x^i, ...x^n\}$          $\{\mathbf{X}^1, \mathbf{X}^2, …\mathbf{X}^i, …\mathbf{X}^n\}$

**Item:**     $\text{Size}(x^i)$:          Sum,
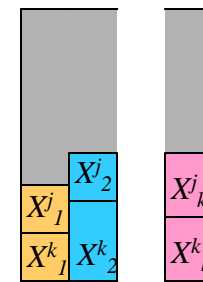          Max, $(X^i_1, X^i_2, …, X^i_k)$:
          Lex.

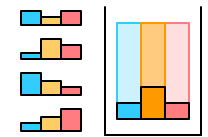**Bin:**     Capacity $C$:          $\mathbf{C} = (C_1, C_2, …, C_k)$:

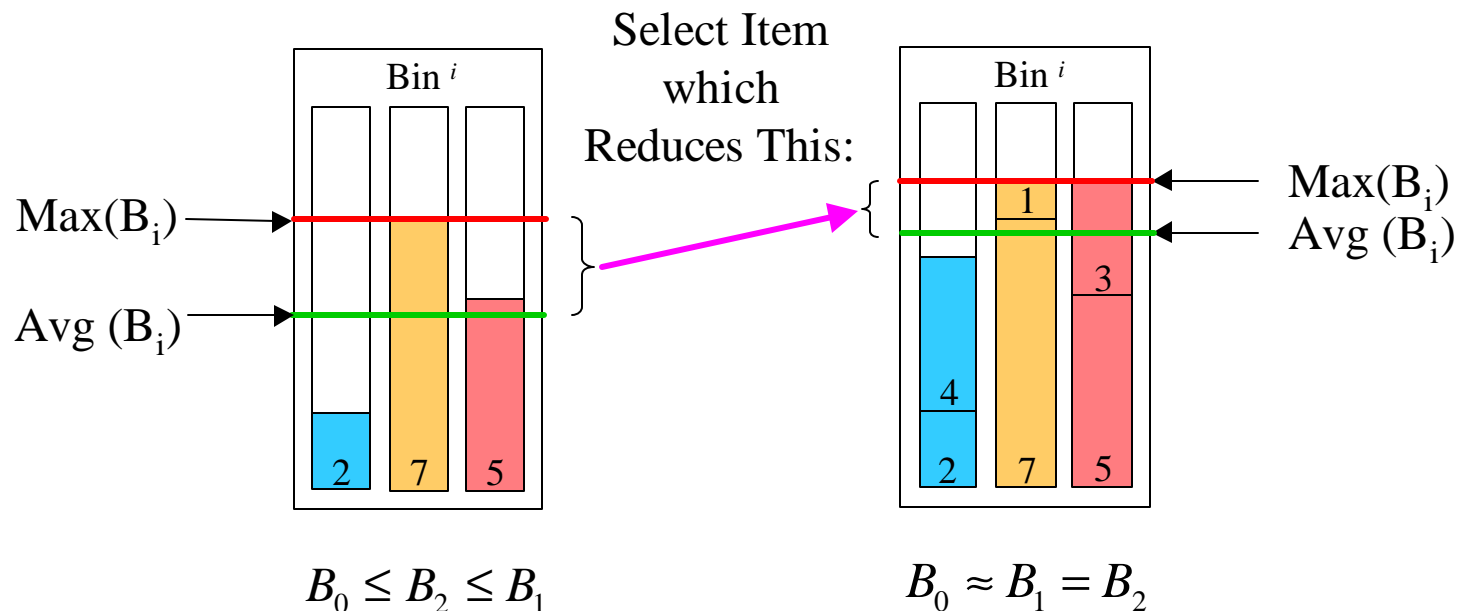          Content $B^j$:          $\mathbf{B}^j = (B^j_1, B^j_2, …, B^j_k)$:

$\Rightarrow$ Goal: Minimize the number of bins required to hold all input items

# K-Capacity Bin Packing: The Resource Balancing Heuristic

Bin Packing

$X^{(i+4)}$ : 3 2 2

$X^{(i+3)}$ : 4 2 4

$X^{(i+2)}$ : 4 1 3

$X^{(i+1)}$ : 4 3 1

$X^i$ : 2 4 5

L

If Bin state is this:

Bin State becomes:

Select Item which Reduces This:

Bin $i$

$Max(B_i)$ →
$Avg\ (B_i)$ →

2 7 5

$B_0 \leq B_2 \leq B_1$

Bin $i$

1
3
4
2 7 5

$Max(B_i)$
$Avg\ (B_i)$

$B_0 \approx B_1 = B_2$

Note:
• $Max(B_i)/Avg(B_i)$ is an approximation to resource balancing
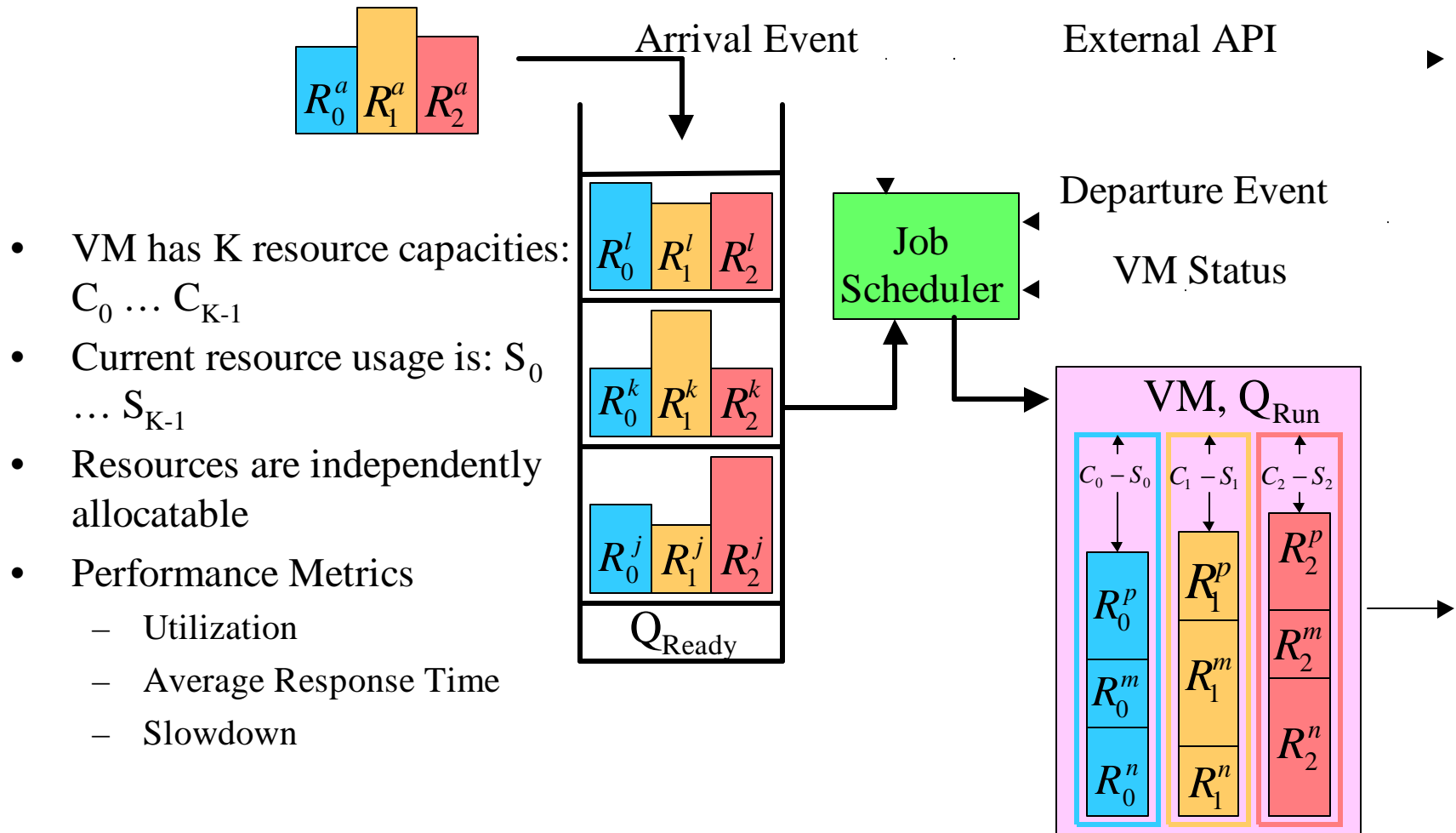
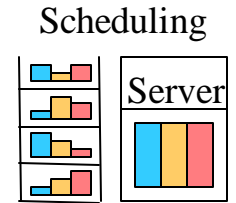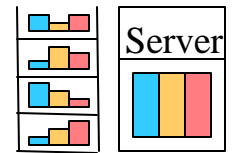HPCMOD UGC 2000

# K-Capacity Bin-Packing: Status

- Progress to Date:
  - Implemented various methods for approximating the Resource Balancing Heuristic.
  - 10-20% reduction in number of bins over First/Best-Fit variants.
  - Results published in ICPP'99.

- Future Work:
  - Item list pre-processing
    - Sorting by size
    - Combining complementary items
  - Dealing with large items and small bins

# K-Resource System Scheduling: Notations

$R_0^a$ $R_1^a$ $R_2^a$

Arrival Event

External API

- VM has K resource capacities: $C_0 \ldots C_{K-1}$
- Current resource usage is: $S_0 \ldots S_{K-1}$
- Resources are independently allocatable
- Performance Metrics
  - Utilization
  - Average Response Time
  - Slowdown

$R_0^l$ $R_1^l$ $R_2^l$

$R_0^k$ $R_1^k$ $R_2^k$

$R_0^j$ $R_1^j$ $R_2^j$

$Q_{Ready}$

Job Scheduler

Departure Event

VM Status

VM, $Q_{Run}$

$C_0 - S_0$ $C_1 - S_1$ $C_2 - S_2$

$R_2^p$

$R_1^p$

$R_0^p$

$R_1^m$

$R_0^m$

$R_2^m$

$R_1^n$
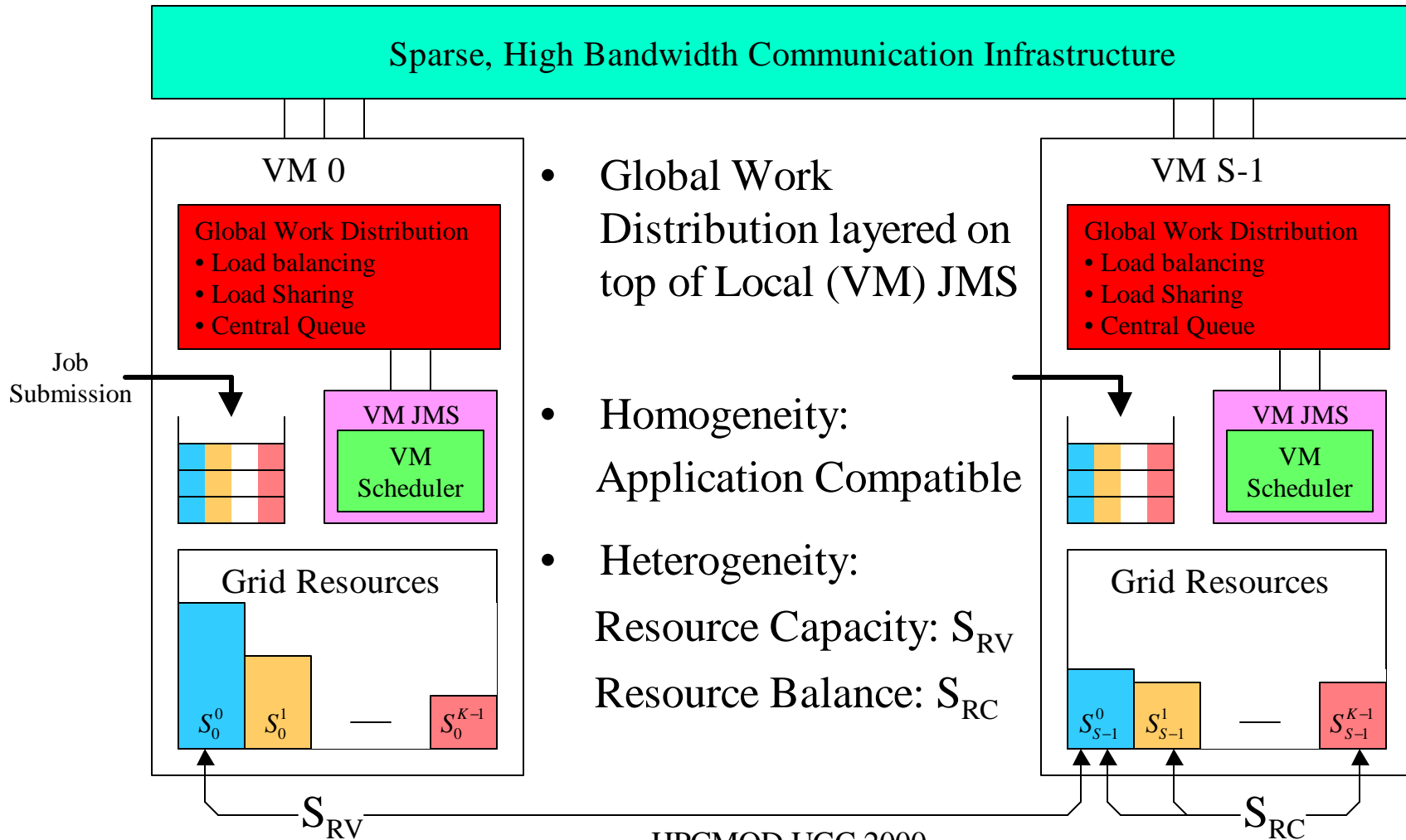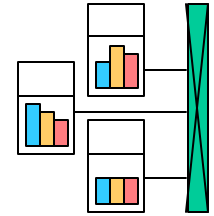
$R_0^n$

$R_2^n$

HPCMOD UGC 2000
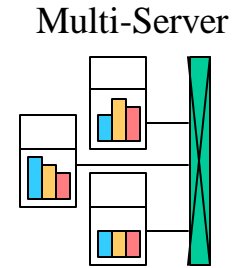
# K-Resource Server Scheduling: Status

- Progress to Date:
  - Resource balancing heuristics are used to select jobs which best match remaining VM resources.
  - Designed new K-Resource aware backfill job selection algorithms.
  - 5-10% gains in Resource Utilization over EASY with First/Best-Fit, on synthetic k-resource workload.
  - 25-50% gains in Average Wait Time over EASY with First/Best -Fit.
  - Results published in SC'99.

- Work in Progress:
  - Unrestricted job queue re-ordering
  - Interaction between the K-Resource grid scheduler and the K single resource local schedulers.

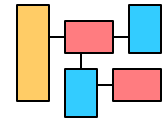# Load Balancing Across Near-Homogeneous K-Resource VMs: Notations

Multi-Server

**Sparse, High Bandwidth Communication Infrastructure**

### VM 0

**Global Work Distribution**
- Load balancing
- Load Sharing
- Central Queue

Job Submission

**VM JMS**

VM Scheduler

**Grid Resources**

$S_0^0$ $S_0^1$ — $S_0^{K-1}$

### VM S-1

**Global Work Distribution**
- Load balancing
- Load Sharing
- Central Queue

**VM JMS**

VM Scheduler

**Grid Resources**

$S_{S-1}^0$ $S_{S-1}^1$ — $S_{S-1}^{K-1}$

- Global Work Distribution layered on top of Local (VM) JMS

- Homogeneity: Application Compatible

- Heterogeneity:
  Resource Capacity: $S_{RV}$
  Resource Balance: $S_{RC}$

$S_{RV}$

$S_{RC}$

HPCMOD UGC 2000

# Load Balancing Across Near-Homogeneous K-Resource Servers: Notations

Multi-Server

- **Progress to Date:**

  - Designed new load balancing policies, based on the resource balancing heuristic, which match the K-Resource capabilities of a server to the K-Resource requirements of the local workload.

  - Makes local scheduler (also K-Resource enabled) more efficient.

  - Achieved performance gains of 5-15% in throughput over classical load balancing techniques.

  - Results published in HCW 2000.

- **Work in Progress:**

  - Load balancing performance is consistently less than central queue.

  - Investigate compromise between central queue and pure load-balancing approaches.

HPCMOD UGC 2000

# Dynamic Creation of a K-Resource
# Virtual Machine in a Computational Grid

Multi-Resource

- How do you select resources?
  - Match the VM resource capabilities to the workload capabilities (Resource Balancing heuristic, again)

- Who's responsible?
  - Grid Infrastructure?
    - Cluster grid applications by resource requirement from a central grid "Ready Queue".
    - Create a K-Resource VM for each cluster/application class.
  - Problem Solving Environments?
    - PSEs used as interface between scientist and grid.
    - Typically multiple instances in use, generating application class specific workloads.
    - PSEs cooperate to make use of grid infrastructure co-allocation services when workload is sufficient to warrant a VM instantiation.
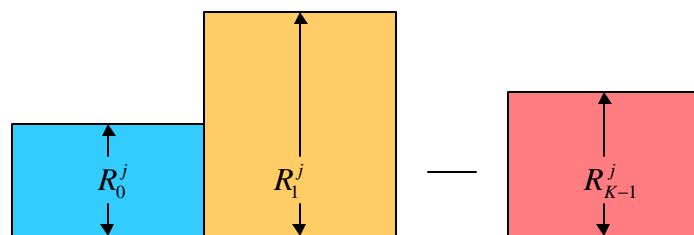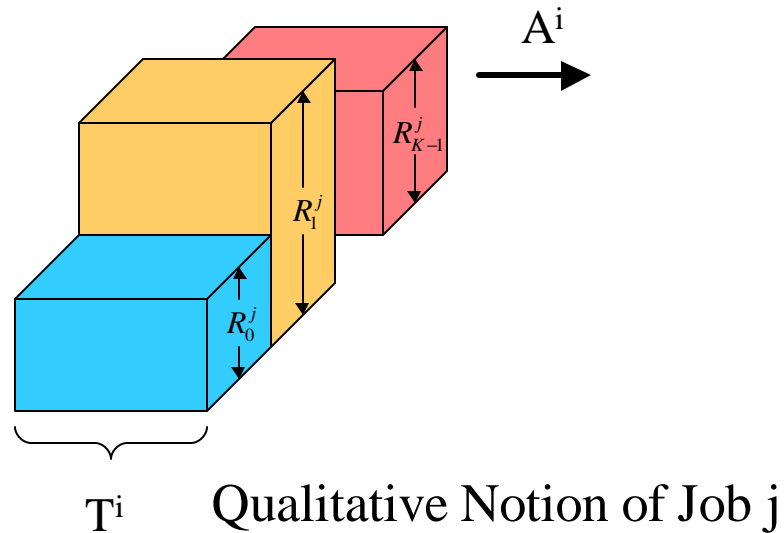
# Summary

- We have proposed the use of K-Resource Virtual Machines to simplify the system-centric scheduling of multi-resource grid-scale applications.

- Our research is incrementally evolving a suite of K-Resource aware scheduling algorithms.

- The mechanisms for creating K-Resource Virtual Machines are extensions of current grid infrastructure services:
  - GLOBUS: A VM is a form of a co-allocation request.
  - LEGION: A VM is a composite object, created from the individual resource objects plus VM scheduler object.

- Still need more data on grid application workloads and system configurations to validate approach.

# Backup Slides

# K-Resource Workload: Notations

$$A^i$$

$$R_{K-1}^j$$

$$R_1^j$$

$$R_0^j$$

$$T^i$$   Qualitative Notion of Job j

$$R_0^j \qquad R_1^j \qquad \text---\qquad R_{K-1}^j$$

K-Resource Job Icon

- Job Assumptions:
  - Each job, j, has K critical resource requirements: $R_0^j \dots R_{K-1}^j$
  - All K resources are required simultaneously
  - The execution time is known, $T^i$
  - Jobs arrive at a rate $A^i$ to a VM system $S^i$
  - Jobs are independent

HPCMOD UGC 2000

# The K-Resource Workload Model: Status

- Progress to Date:
  - Created a synthetic model for a K-Resource workload
    - Simple extension of single resource (CPUs) parallel processing workload models
    - Inter-Job resource probability distribution
    - Intra-Job resource correlation

- Work in Progress:
  - Searching for characteristics of grid-scale application workloads
    - Resource configuration requirements for application classes, e.g. collaborative design environments, real-time experiment control, distributed parallel processing, etc..
    - Characteristics of resource classes, e.g. data source/sink, network link, compute engine, visualization, storage devices, etc..

HPCMOD UGC 2000

# Relevant Research

- Job Classifications
  - Rigid, Evolvable, Modable, Malleable
- Workload Trace Analysis
  - General (size=CPUs)
    - Size distributions are exponential/hyperexponential
    - Large percentage from discrete sizes - powers of 2, squares of integers
    - Execution time is weakly correlated to job size
  - Memory
    - Many discrete sizes as well
    - Weakly correlated to number of CPUs
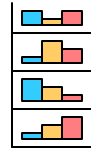  - I/O
    - Input, Computation, Output Phases
    - Bursty

# Initial Approach:
# Inter-Job Resource Distributions

Workload

$$D_{\beta(m,v)}, D_L, p_L$$



Large:
- powers of 2
- integer squares

$$D_L^{Range}$$

Standard Beta
- mean      0.10
- variance 0.09

$$D_{\beta(m,v)}^{Range}$$

$p_L = 0.05$

$p_L = 0.10$

$p_L = 0.15$

$p_L$

HPCMOD UGC 2000

# Initial Approach:
# Intra-Job Resource Correlation, $R_C$

$K=2$, $p_L=0.1$, $\beta(0.10, 0.09)$

HPCMOD UGC 2000

# K-Resource VM Scheduling: Progress
## The Resource Balancing Heuristic

Scheduling

Server

3 2 2

$X^{(i+4)}$

4 2 4

$X^{(i+3)}$

4 1 3

$X^{(i+2)}$

4 3 1

$X^{(i+1)}$

2 4 5

$X^i$

L

If VM state is this:

VM State becomes:

Select Job which Reduces This:

Server

$\text{Max}(S_i) \longrightarrow$

$\text{Avg}(S_i) \longrightarrow$

2 7 5

$S_0 \leq S_2 \leq S_1$

Server

1

3

4

2 7 5

$\longleftarrow \text{Max}(S_i)$

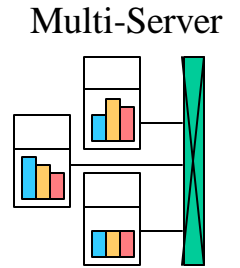$\longleftarrow \text{Avg}(S_i)$

$S_0 \approx S_1 = S_2$

Note:
- $\text{Max}(S_i)/\text{Avg}(S_i)$ is an approximation to resource balancing

HPCMOD UGC 2000

# Components of Classical Load Balancing

| Policy | Description | Baseline Approach |
| --- | --- | --- |
| **Load Index** | Describes Load on each Server | Job Queue Lengths, Resource Average (RA) |
| **Information Policy** | How and When to update Load Index | Global, Instantaneous |
| **Transfer Policy** | Triggering a load balancing action | Threshold-Initiated Sender, Receiver, and Symmetric |
| **Selection Policy** | Which job to move | Latest Job Arrived (LSP) |
| **Location Policy** | Where to move job from/too | Lowest Load, Highest Load |

HPCMOD UGC 2000

# Limitations of Resource Average (RA) Load Index in the Presence of Multiple Resources

Multi-Server



Q⁰ — $Q^0$

Q¹ — $Q^1$

QG — $Q^G$

Load Index

RB⁰ — $RB^0$

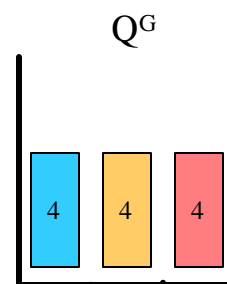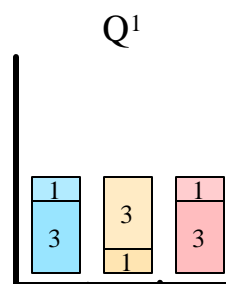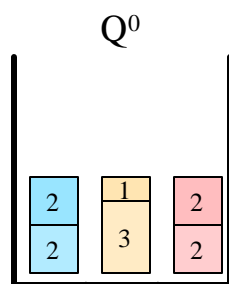RB¹ — $RB^1$

Resource Average (RA)

Resource Balance (RB)

- The relative total resource requirements of the Workload should match the relative resource capabilities of the Server.
  - e.g. The workload of a server with a large memory configuration should contain mostly memory intensive jobs.
- Extend the classical Load Index, RA to RA+RB, to trigger load balance actions in the event of a workload resource balance mismatch.

HPCMOD UGC 2000

# Limitations of Latest Job (BSP) Selection Policy in the Presence of Multiple Resources

Multi-Server

$Q^0$  $Q^1$  $Q^G$

$RB^0$

$RB^1$

**Resource Average (RA)**

**Resource Balance (RB)**

$Q^0$  $Q^1$  $Q^G$

**Last Job Selection Policy (LSP)**

- Doesn't actively correct workload balance

$Q^0$  $Q^1$  $Q^G$

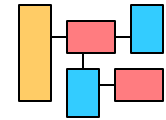**Balanced Selection Policy (BSP)**

- Actively corrects workload balance

- Add new job selection policy, BSP, to actively correct a workload resource balance mismatch.

# Dynamic Creation of K-Resource Virtual Machines in a Computational Grid

- ## When Should a K-Resource VM be Created?
  - Upon submission of a single grid application? No!
    - Leads to First-Come-First-Served use of critical resources
  - Problem Solving Environments? Yes!
    - PSEs used as interface between scientist and grid
    - Typically multiple instances in use, generating application class specific workloads
    - PSEs cooperate to make calls to grid infrastructure VM services when workload is sufficient

- ## Creating a K-Resource VM in the Grid Infrastructure:
  - Grid Infrastructure
    - GLOBUS: Add service to create VM, similar to co-allocation request.
    - LEGION: Object to create a VM composite object.

HPCMOD UGC 2000